

(12) INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(19) World Intellectual Property Organization
International Bureau



(43) International Publication Date
22 May 2003 (22.05.2003)

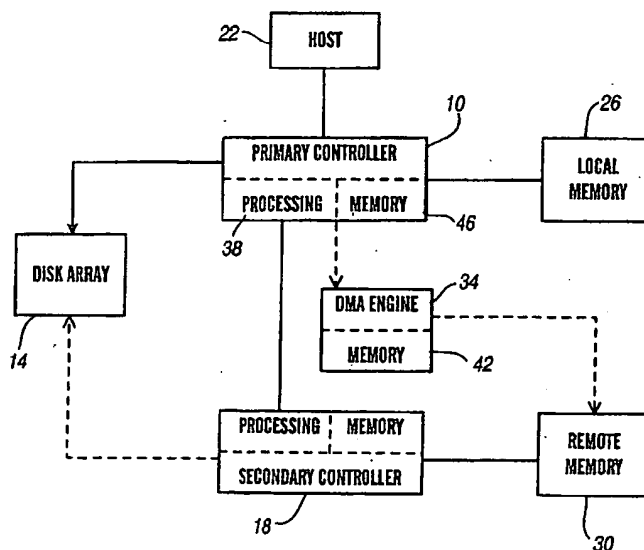
PCT

(10) International Publication Number
WO 03/043254 A2

- (51) International Patent Classification⁷: **H04L** (74) Agents: **ZINGER, David, F. et al.**; Sheridan Ross P.C., Suite 1200, 1560 Broadway, Denver, CO 80202-5141 (US).
- (21) International Application Number: **PCT/US02/35786** (81) Designated States (*national*): AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DZ, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NO, NZ, OM, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, YU, ZA, ZM, ZW.
- (22) International Filing Date:
7 November 2002 (07.11.2002)
- (25) Filing Language: English
- (26) Publication Language: English
- (30) Priority Data:
60/332,415 9 November 2001 (09.11.2001) US
- (71) Applicant (*for all designated States except US*): **CHAPARRAL NETWORK STORAGE, INC.** [US/US]; 7420 East Dry Creek Parkway, Longmont, CO 80503 (US).
- (72) Inventor; and (84) Designated States (*regional*): ARIPO patent (GH, GM, KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZM, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE, SK, TR), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG).
- (75) Inventor/Applicant (*for US only*): **MAINE, Gene** [US/US]; 1622 Albion Lane, Longmont, CO 80503 (US).
- Published:
— *without international search report and to be republished upon receipt of that report*

[Continued on next page]

(54) Title: TRANSFERRING DATA USING DIRECT MEMORY ACCESS



(57) Abstract: A direct memory access (DMA) engine has virtually all control in connection with data transfers that can involve one or both of primary and secondary controllers. The DMA engine receives a command related to a data transfer from a processor associated with the primary controller. This command causes the DMA engine to access processor memory to obtain metadata therefrom. In performing a DMA operation, the metadata enables the DMA engine to conduct data transfers between local memory and remote memory. In performing exclusive OR operations, the DMA engine is involved with conducting data transfers using local memory.

WO 03/043254 A2



For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.

TRANSFERRING DATA USING DIRECT MEMORY ACCESS

FIELD OF THE INVENTION

The present invention relates to direct memory access engines, and more particularly,
5 to a method and apparatus for transferring data between RAID controllers using a direct memory access engine.

BACKGROUND OF THE INVENTION

Redundant Arrays Of Inexpensive Disks (RAID) systems are well known systems
10 which can help increase availability of stored data in network storage systems. Such systems typically include several hard disk drives which store data in such a way that, if one disk drive fails, data is still able to be recovered from the system. Such systems generally have a network storage bridge, which acts as an interface between a host computer and an array of disk drives. To further enhance the availability of data in such a system, it is common to
15 have redundant storage controllers within the network storage bridge, such that if one controller fails, the remaining controller is able to continue read and write operations to the array of disk drives.

In a fully redundant storage bridge, user data has to be temporarily stored twice, once in a primary controller of data and once in its redundant or secondary counterpart. In this
20 setup if the primary controller is damaged and unable to continue operation, the secondary redundant controller has a copy of the user data. Thus data is available if one controller fails, because either the primary or the secondary copy of the data will be delivered to its final destination, the array of disk drives.

Typically and with reference to the prior art drawing of Fig. 1, when a primary
25 controller 10 receives user data to be stored in the disk array 14, it builds description tables in the process. These tables describe the data size, location and its destination. The primary controller 10 then initiates a data transfer, known as a mirroring operation, in which the data is transferred to the secondary controller 18. Once the mirroring operation is complete, the primary controller 10 sends an indication to the host computer 22 that the data write is
30 complete. Thus, the host computer 22 is notified that the data is stored, while the data may not be written to the disk array 14 for a period of time.

A traditional method for mirroring data between controllers is for the primary controller 10 to notify the secondary controller 18 that data is going to be mirrored. The

primary controller 10 then transfers metadata to the secondary controller 18. The metadata is the data which is contained in the description tables. Following the metadata, the primary controller 10 transfers the user data to the secondary controller 18. The user data is stored in a local memory 26 associated with the primary controller 10 for transfer to a remote
5 memory 30 associated with the secondary controller 18. In some cases, these transfers are initiated using direct memory access (DMA) and a DMA engine 34. In such a situation, a processing portion 38 within the primary controller 10 will give the DMA engine 34 a transfer command to transfer data from the primary controller 10 to the secondary controller 18. This transfer command typically includes the data contained in the description tables,
10 which identifies the data which is to be transferred to the secondary controller 18. The processing portion 38 generally builds a DMA table containing the data from the description tables, which is able to be used by the DMA engine 34. This DMA table is then loaded into the DMA engine 34, with the transfer command. The DMA engine 34 receives the data from the processing portion 38, stores it in a memory 42 associated with the DMA engine 34, and
15 then conducts the data transfer.

As can be seen, the processing portion 38 must build a DMA table, and transfer the table from the processing portion to the DMA engine 34. Thus, this data is stored in a memory 46 associated with the processing portion 38, configured into a form which is usable by the DMA engine 34, and then transferred and stored in the memory 42 associated with the
20 DMA engine 34. It would be advantageous to reduce the amount of processing overhead involved in a mirroring transaction, thereby improving system performance. Accordingly, it would be advantageous to have a DMA engine which does not require a processing portion to create a DMA table. Furthermore, it would be advantageous to reduce internal memory required in a DMA engine, thereby reducing the silicon area required for such a DMA
25 engine.

SUMMARY OF THE INVENTION

This invention allows a hardware based DMA engine to use the description tables stored in the processing portion directly to transfer data from the primary controller to the
30 secondary controller through a dedicated, redundant set of data paths internal to the data bridge subsystem. Thus, the processing portion of the primary controller is not required to

build a DMA table for use by the DMA engine. This transfer is accomplished via an internal data path from the primary controller to the secondary controller without ever leaving the storage bridge itself. Furthermore, the DMA engine uses the data tables as-is, without modification, thus reducing processor overhead. As the data is being transferred the tables
5 are being updated and upon completion of a transfer, a completion status is posted to the processing portion. Since the DMA engine accesses the data table present in the processing portion memory, the processing portion does not have to load or build the table into the DMA memory. This access is accomplished independently of the processor portion thereby reducing processor overhead, as well as allowing for a reduced memory size for the DMA
10 engine, saving silicon area for the chip containing the DMA engine.

The engine is also capable of XOR (exclusive OR) operations on the user data based on the same tables as the DMA transfer. In fact the XOR and DMA of the data can be performed concurrently, thus reducing overhead further.

The DMA engine of this invention thus results in a savings in gate count and silicon
15 area required for its implementation, a reduction in CPU (central processing unit) processing time and thus reduction in command overhead, and reduction in software development time, since no additional data structures are needed for the DMA engine, which transfers data independently of the processor.

20 BRIEF DESCRIPTION OF THE DRAWINGS

Fig. 1 is a block diagram illustrating a prior art system that includes primary and secondary controllers which can include and a DMA engine in connection with the storing of data;

Fig. 2 is a block diagram of a controller with a DMA engine for performing DMA
25 and/or XOR operations upon receiving a CPU command;

Fig. 3 is a block diagram of the DMA engine that includes a DMA engine portion and a XOR engine portion;

Fig. 4 illustrates a command queue associated with the DMA engine;

Fig. 5 illustrates a scatter-gather list structure;

30 Fig. 6 illustrates scatter-gather lists associated with a XOR function; and

Fig. 7 is a flow diagram illustrating operations associated with the DMA engine.

DETAILED DESCRIPTION

The invention is implemented as a module inside a memory controller 50. This device, illustrated in Fig. 2, moves data across four ports. Two multi-clients PCI-X ports 54, 58, one of which connects the primary controller 50 and the secondary controller (not shown), a CPU port 62 (PCI Interface Bus C) and a DDR-SDRAM port 66 for temporary and
5 cache data storage. The DMA engine 70 is one of the internal modules of this device.

Fig. 3 illustrates the internal partitioning of the DMA engine 70. The DMA engine 70 is started by a CPU command issued via the CPU port 62. The CPU command includes the location of the data description table, located in the CPU memory 74. The command also
10 includes instructions on the type of command to execute (DMA or XOR) and the PCI-X port 54, 58 to use for the DMA command. Once the command is received the DMA engine 70 reads the CPU memory based tables and performs data movement and manipulation without further CPU 78 assistance. As data is being processed, the DMA engine 70 updates (writes) the CPU tables to indicate how far the operation has progressed. The DMA engine 70 does
15 not store the tables internally. Instead the CPU tables are periodically accessed through the dedicated CPU port 62 as more data is being transferred. This allows for reduced processor overhead, as the CPU 78 is not required to build a DMA table or load DMA transfer information into the DMA engine 70.

The data manipulated by the DMA engine 70 does not have to be contiguous. The
20 CPU tables typically represent the scatter-gather nature of the RAID data storage (see Figs. 3, 4 and 5). These tables tend to describe data in RAID chunks and these are the same tables that the DMA engine 70 uses to copy data across redundant controllers as well as XOR the data inside the primary controller 50 or the secondary controller.

The DMA table consists of destination entries only. Each table entry consists of
25 sub-fields referred to as elements. Each element describes the location and the amount of data to copy. The "LAST" bit at the end of the source and destination tables instructs the DMA engine 70 to complete the operation.

The data copy table consists of source and destination entries. This command is designed to copy data inside the local controller memory 82. Each element describes the
30 location and the amount of data to copy. The source table describes the locations and amounts of data located in the primary controller 50, the destination table describes the

locations in the secondary controller to which the data has to be moved. The "LAST" bit at the end of the source and destination tables instructs the DMA engine 70 to complete the operation.

The XOR table resembles the data copy table, except it has multiple source tables. Data location and amount described by these source tables is XORed and the result is moved to the locations described by the destinations table. Again the "LAST" bit at the end of the source and destination table instructs the DMA engine 70 to complete the operation. In addition the "Source Count" field alerts the DMA engine 70 as to the number of source data pointers to XOR. Since the tables are not stored inside the DMA engine 70 itself, the theoretical number of sources is unlimited and does not impact the size of the DMA engine 70.

Once the DMA engine 70 operation has started, the system CPU 78 has an option to enqueue another command into the DMA engine 70 thus reducing command overhead to almost zero or alternatively to monitor the progress and wait for a command completion interrupt. Once the DMA engine 70 has completed all required data transfer, it issues a maskable interrupt and a status to the CPU 78 and ceases operation.

The DMA engine 70 has additional system benefits. It can be used to quickly initialize all memory to a predetermined data pattern. Another benefit of the DMA engine 70 is that it can be used to quickly check memory for errors, if a checksum or an ECC is used in conjunction with the temporary storage memory, as is the case in typical redundant systems.

Both the XOR engine portion 86 and the DMA engine portion 90 of the DMA engine 70 use common Scatter-Gather list structures. The S-G lists reside in the CPU memory 74 and the DMA engine 70 is responsible for extracting information from these lists and performing necessary data movement and manipulation. Unlike the S-G lists, the source of the actual user data being manipulated is always the local DDR-SDRAM user data buffer. Data transfer destination for XOR operations is also the local DDR-SDRAM memory 82. Data destination for DMA operations is either the PCI-XA port 54 or the PCI-XB port 58, depending on the command.

Fig. 4 is an illustration of the command queue for the DMA engine 70. Enqueueing and monitoring progress of XOR or DMA operations is performed via several command

registers. Some control and status reporting is performed in the interrupt/status registers located in the interrupt control block. The enqueue register is used to issue commands. It is written with the type of command to be executed and its parameters. The current command register is a read only register which contains statuses of commands currently being executed.

- 5 The status register is used to verify and acknowledge successful or unsuccessful command completion.

Within the command registers are several fields. The memory commands start immediately after the command (CMD) field of the enqueue register. Currently defined legal commands are:

- 10 01 - XOR
 02 - DMA through PCI-X A bus
 03 - DMA through PCI-X B bus

All other values in this field are reserved.

- 15 The source count field (SrcCnt) defines the number of source scatter-gather (S-G) lists that the command is required to process. There is only one destination S-G list. The source count field is only relevant for XOR commands. The DMA commands only use the one allowed destination list as both source and destination, and do not check the SrcCnt field. The current command register displays the current (live) source list number being processed. The S-G address pointer field is used to identify which S-G list structure is to be used for a
- 20 given command. The field is presumably used as both a starting address of the list (when the command is enqueued) as well as a command identifier (when status register is examined). The current command register will always show a "live" version of the address pointing to the list currently being processed. Note that this address pointer is required to point to an existing address in the CPU memory range. The stat field reflects command initialization
- 25 (start) and completion statuses. A value of "00" in the status register signifies that there are no statuses pending. Clearing the status field will automatically make the register available for the initialization/completion status of the next command. Note that clearing this bit also clears the interrupt associated with this command's initialization or completion, provided the interrupt is enabled.

- 30 The command progress register reflects the flow of the commands being enqueued, currently executed and completed. As commands are being enqueued, a CMD_ENQD bit in

this register gets asserted to signify that there is a command that has been enqueued but has not yet started. This bit will continue to be asserted until the command starts. If no other command is being executed, this will occur immediately, however if there is another command in the queue being executed, the bit will stay ON until the previous command execution is completed. Another command can be enqueued only after the currently
5 enqueued command has been started and this bit has been cleared by the hardware.

Once the command starts, the CMD_INPROG bit is asserted. At the same time a maskable "Command Started" interrupt is set. The significance of this event is that once this interrupt is set, another command can be enqueued. Again, this new command will stay
10 enqueued until the current command execution is completed.

Upon completion of either the DMA or the XOR operation, the CMD_COMPL bit is asserted signifying that the command has completed all data movement and/or manipulation. If the command status register is empty (no previous command completion statuses are pending), assertion of this bit will be momentary. However, if the status register
15 is busy, this bit will stay on signifying that the command is waiting to post its completion status as soon as the status register is available for posting.

The command status register contains a valid status when a maskable command completion interrupt is set. The CMD_STAT bit in the command progress register will also be set at the same time. Clearing status bits in the command status register clears both the
20 interrupt and the CMD_STAT bit. It also makes the status register available for the next command completion status.

The S-G list structure is illustrated in Fig. 5. Each individual list consists of elements. Each element represents a starting address of a block of user data and the length of this data. The starting address can be anywhere in the 4GB space. The length of a single
25 element is currently limited to a maximum of 64 Kbytes size, although the DMA engine 70 can be reprogrammed to accept larger sizes. The S-G LAST element flag determines the number of elements in a S-G list. The flag is asserted by placing a '1' in bit 31 of the length field of the last S-G element.

Fig. 6 illustrates an XOR function capable of being completed using the DMA engine
30 70. The entire S-G list structure as used by the XOR operation follows. Only the first destination list is used for DMA operations. In this case it becomes both the destination and

the source list. The source and the destination addresses may optionally differ by 1GB (as defined by the DMA offset enable register). Destination address conversion is handled in hardware. The XOR engine portion 86 can also handle a unique case of a single source list. Although this operation is still initiated as an XOR command, in reality it performs a copy of user data from the source to the destination location. No actual XOR operation is performed. The data is moved from the source location in local DDR-SDRAM to a destination location in DDR-SDRAM memory. The data is not moved across either of the PCI-X buses.

Fig. 7 is a flow diagram illustration of the operation of a DMA engine 70 of one embodiment of the present invention. Initially, as noted by block 100, the CPU 78 initiates a mirroring operation. At block 104, the CPU 78 places a DMA command in the DMA engine 70 command queue. The DMA engine 70, when the DMA command is in the current command register, reads the current command, as noted by block 108. The DMA engine 70 determines if the command is an XOR command, as noted by block 112. If the command is not an XOR command, at block 116, the DMA engine portion 90 retrieves the first scatter/gather element from the CPU memory address indicated by the DMA command. The DMA engine portion 90 copies the data from the local DDR-SDRAM memory 82 location indicated by the first scatter/gather element to a remote DDR-SDRAM (not shown) location which corresponds to the location indicated by the first scatter/gather element, as noted by block 120.

At block 124, the DMA engine portion 90 determines if the last element flag is set in the scatter/gather element. If the last element flag is set, the DMA engine portion 90 marks the command status as complete in the command queue and sends a command complete message to the CPU 78, as noted by block 128. The mirroring operation is then complete, as noted by block 132. If at block 124 the DMA engine portion 90 determines that the last element flag is not set, the DMA engine portion 90 retrieves the next sequential scatter/gather element from CPU memory 74, as noted by block 136. The DMA engine portion 90, according to block 140, then copies data from the local DDR-SDRAM memory 82 location indicated by the scatter/gather element to a remote DDR-SDRAM location which corresponds to the location indicated by the scatter/gather element. The DMA engine portion 90 then repeats the operations associated with blocks 124 through 140.

If, at block 112, the DMA engine 70 determines that the DMA command is an XOR command, the XOR engine portion 86 of the DMA engine 70 reads the source count field from the DMA command register, as noted by block 144. The XOR engine portion 86 then, at block 148, retrieves the first scatter/gather element from the CPU memory address indicated in the DMA command. The XOR engine portion 86 performs XOR operations on the data contained in the memory locations which are determined from the memory location indicated by the first scatter/gather element and stores the XOR result in a destination memory location which is indicated by the memory location in the first scatter/gather element, as indicated by block 152. The XOR engine portion 86 then determines, at block 156, whether the last element flag is set in the scatter/gather element. If the last element flag is set, the XOR engine portion 86, at block 160, marks the command as complete in the command register and sends a notification to the CPU 78 that the command is complete. The XOR operation is then complete, as noted by block 164. If, at block 156, the XOR engine portion 86 determines that the last element flag is not set in the scatter/gather element, the XOR engine portion 86, at block 168, retrieves the next sequential scatter/gather element from CPU memory 74. The XOR engine portion 86, at block 171, performs XOR operations on the data contained in the memory locations which are determined from the memory location indicated by the scatter/gather element and stores the XOR result in a destination memory location which is indicated by the memory location in the scatter/gather element. The XOR engine portion 86 then repeats the operations associated with blocks 156 through 172.

The foregoing discussion of the invention has been presented for purposes of illustration and description. The description is not intended to limit the invention to the form disclosed herein. Variations and modifications commensurate with the above teachings, within the skill and knowledge of the relevant art, are within the scope of the present invention. The embodiment described hereinabove is further intended to explain the best mode presently known of practicing the invention and to enable others skilled in the art to utilize the invention in such embodiment, or in other embodiments, and with the various modifications required by their particular application or use of the invention. It is intended that the appended claims be construed to include alternative embodiments to the extent permitted by the prior art.

What is claimed is:

1. A method for processing data, comprising:
providing at least a first controller communicable with a remote memory, said first controller including a direct memory access (DMA) engine;
5 providing at least a processor communicating with a processor memory;
storing information using said processor with said processor memory related to data for processing, said information comprising a number of elements including at least a first element and a second element;
starting a process related to said information;
10 accessing firstly said processor memory using said DMA engine to obtain said first element;
processing first data related to said first element;
accessing secondly said second element from said processor memory using said DMA engine after said step of processing said first data; and
15 processing second data related to said second element.
2. A method, as claimed in Claim 1, wherein:
said step of accessing secondly is conducted independently of said processor.
3. A method, as claimed in Claim 1, wherein:
said accessing firstly is conducted with a data format for said first element that is the
20 same data format of said first element that was used when it was stored in said processor memory.
4. A method, as claimed in Claim 1, wherein:
at least said first controller is part of a first structure including a chip and said processor memory is external of said structure.
- 25 5. A method, as claimed in Claim 1, wherein:
said starting step includes sending a command to said DMA engine.
6. A method, as claimed in Claim 1, wherein:
for each of said number of elements, said DMA engine accesses said processor memory without interrupting said processor and, for each of said accessing, a processing of
30 data is conducted before another accessing is conducted of said processor memory using said DMA engine.

7. A method, as claimed in Claim 1, wherein:
said processing first data includes copying said first data to said remote memory that
is relatively remote from said DMA engine.
8. A method, as claimed in Claim 1, wherein:
5 said processing first data includes performing a XOR operation.
9. A method for mirroring data from a local memory to a remote memory,
comprising:
receiving a command at a direct memory access (DMA) engine;
retrieving by said direct memory access engine information from a memory location
10 external to said direct memory access engine;
copying, based on said information contained in said memory location, user data from
a local memory location to a remote memory location.
10. A method, as claimed in Claim 9, wherein:
said command is a CPU command related to at least one of a DMA operation and a
15 XOR operation.
11. A method, as claimed in Claim 9, further including:
determining using said DMA engine whether a DMA operation or a XOR operation
is to be performed.
12. A method, as claimed in Claim 9, wherein:
20 said information includes description data related to at least a location of said user
data.
13. A method, as claimed in Claim 9, wherein:
said memory location communicates with a CPU.
14. A method, as claimed in Claim 9, wherein:
25 said information includes a number of scatter-gather elements.
15. A method, as claimed in Claim 9, wherein:
said command is sent by a CPU and said copying step is conducted substantially
independently of said CPU.

16. An apparatus for mirroring data in a storage system, comprising:
a local memory operable to store user data;
a remote memory operable to store user data;
a processor including a processor memory; and
5 a direct memory access (DMA) engine operable to retrieve description data from said processor memory and to copy said user data from said local memory to said remote memory based on said description data.
17. An apparatus, as claimed in Claim 16, wherein:
said DMA engine receives a command related to retrieving said description data and
10 copying said user data from said processor and in which said command includes a location in processor memory related to said description data.
18. An apparatus, as claimed in Claim 16, wherein:
said DMA engine includes at least one command register that has a plurality of fields including a command field related to at least one of an exclusive OR operation and a DMA
15 operation and a source count field related to a number of scatter-gather lists.
19. An apparatus, as claimed in Claim 18, wherein:
said at least one command register also includes a field related to the current source list number being processed and a field related to an address pointer that is used to identify at least one of said scatter-gather lists associated with said command.
20. An apparatus, as claimed in Claim 16, wherein:
said DMA engine includes a DMA engine portion used in performing DMA operations and a XOR engine portion used in performing exclusive OR operations.

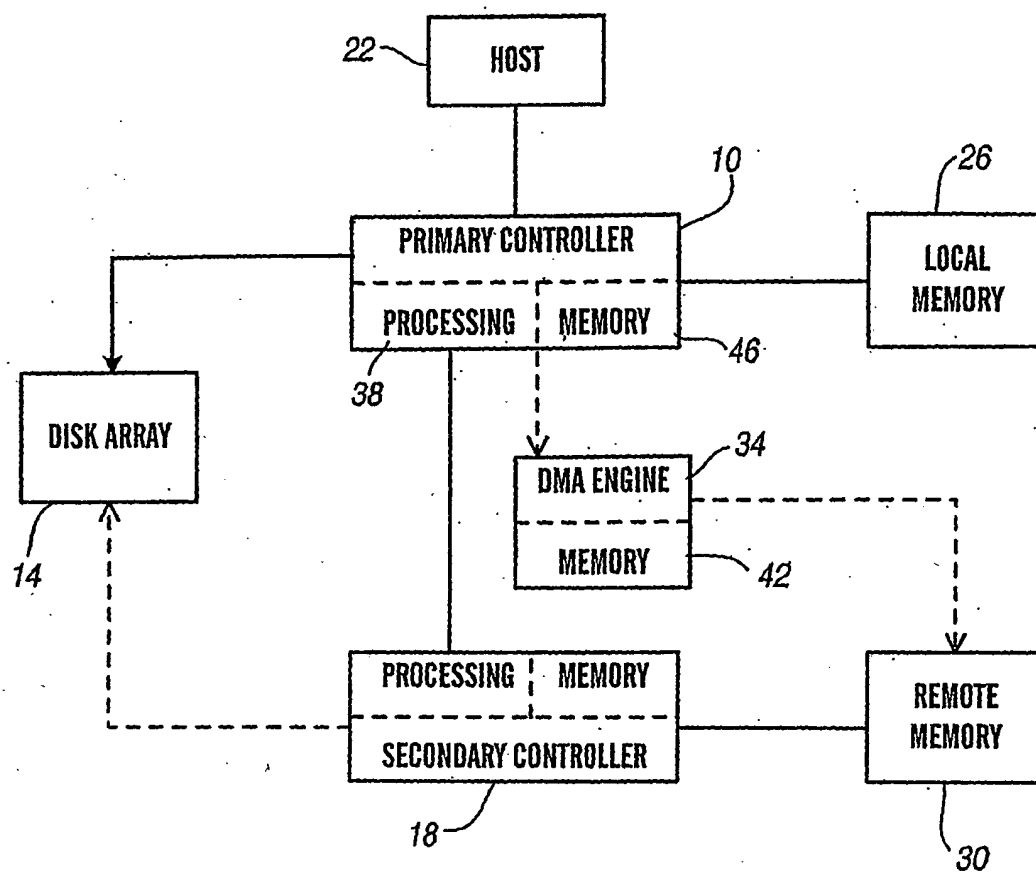


FIG. 1
(PRIOR ART)

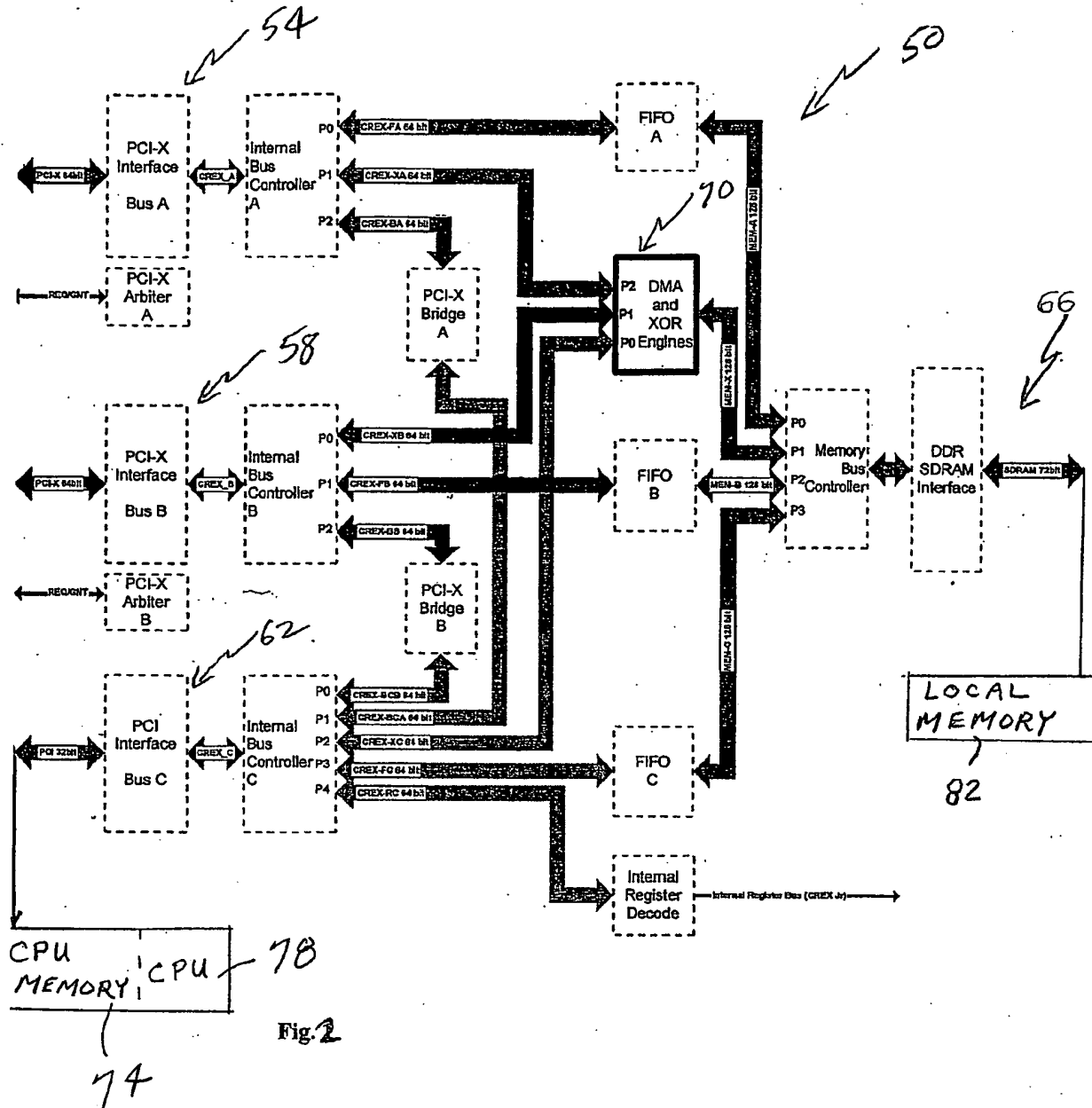


Fig. 2

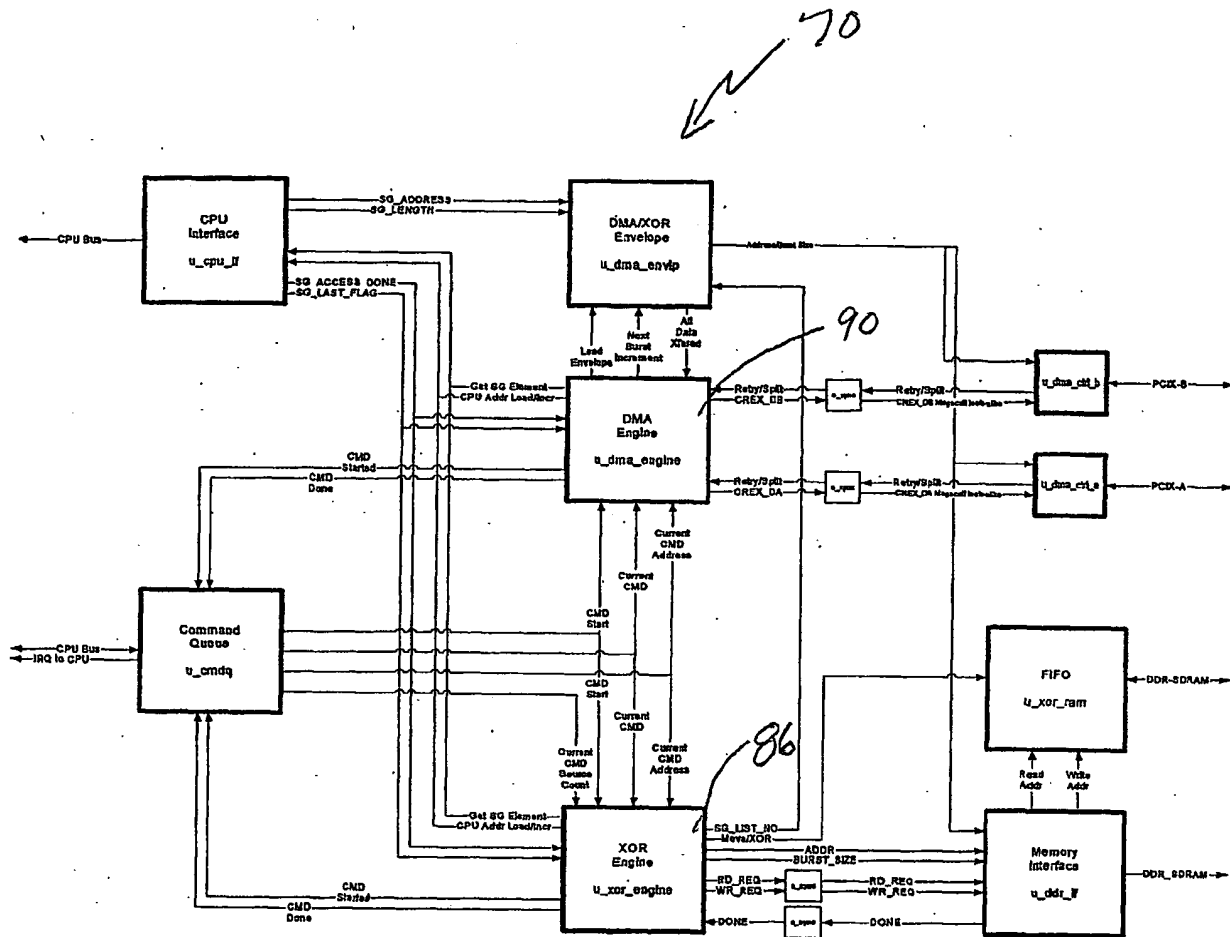


Fig. 3

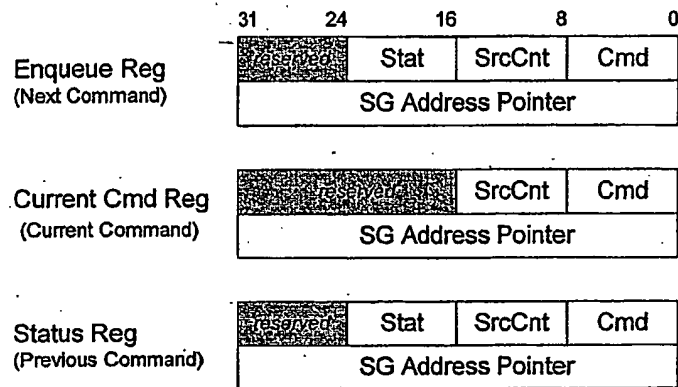


Fig. 4

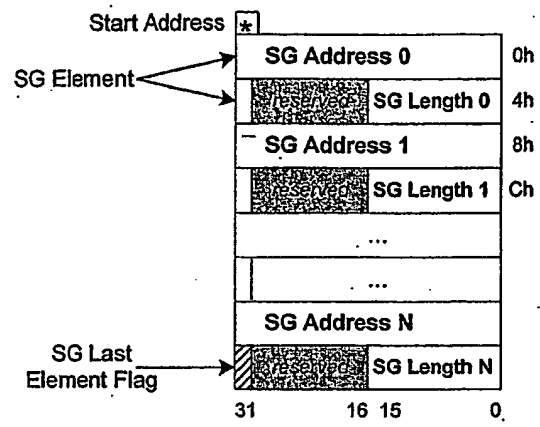


Fig. 5

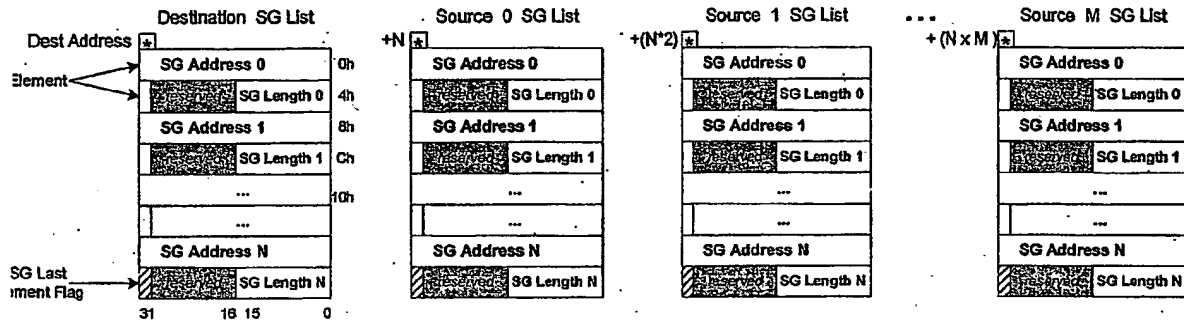


Fig. 6

